

Parsed Page eXplorer (PPX): A Semantic Search Engine for Document Understanding

Levi Willms • Ken Pu

Ontario Tech University • Oshawa, ON, Canada

levi.willms@ontariotechu.net, ken.pu@ontariotechu.ca

1. Motivation

- **Problem:** Raw PDF documents lack machine-readable structure, making semantic queries impossible to spatially locate. Existing OCR tools provide bounding boxes but fragment text; LLMs reconstruct meaning but lose spatial information.
- **Challenge:** Aligning spatial OCR fragments with semantic LLM text to enable queries that return precise page coordinates.
- **Goal:** Establish semantic queries on PDF documents with verifiable source locations — establishing *semantic provenance* (tracing answers to their exact spatial origin).

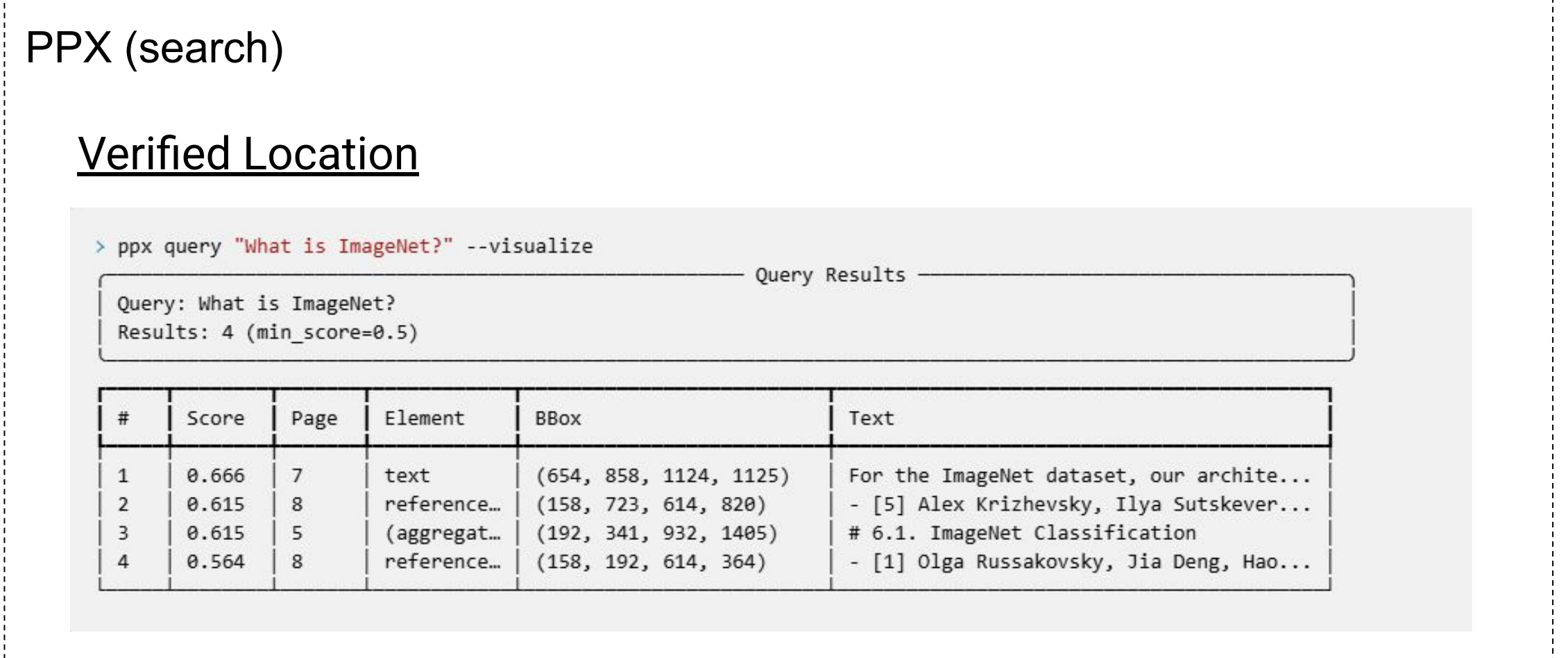
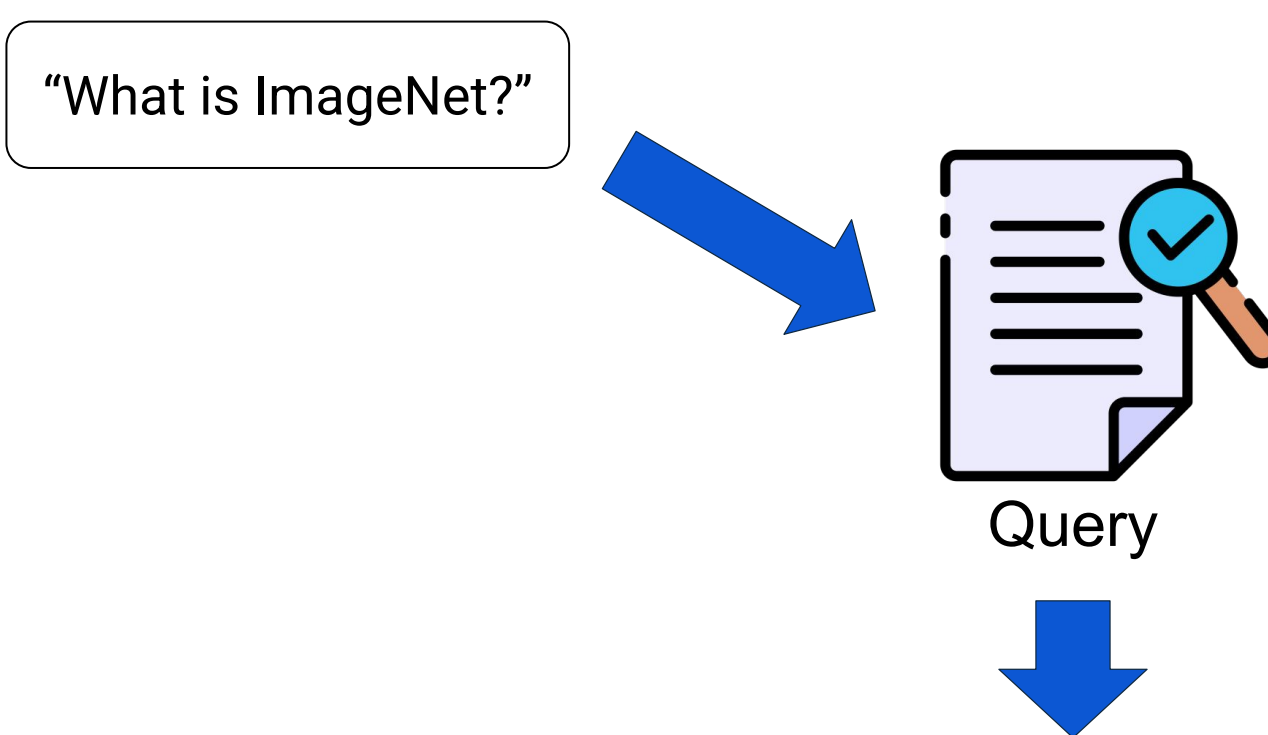
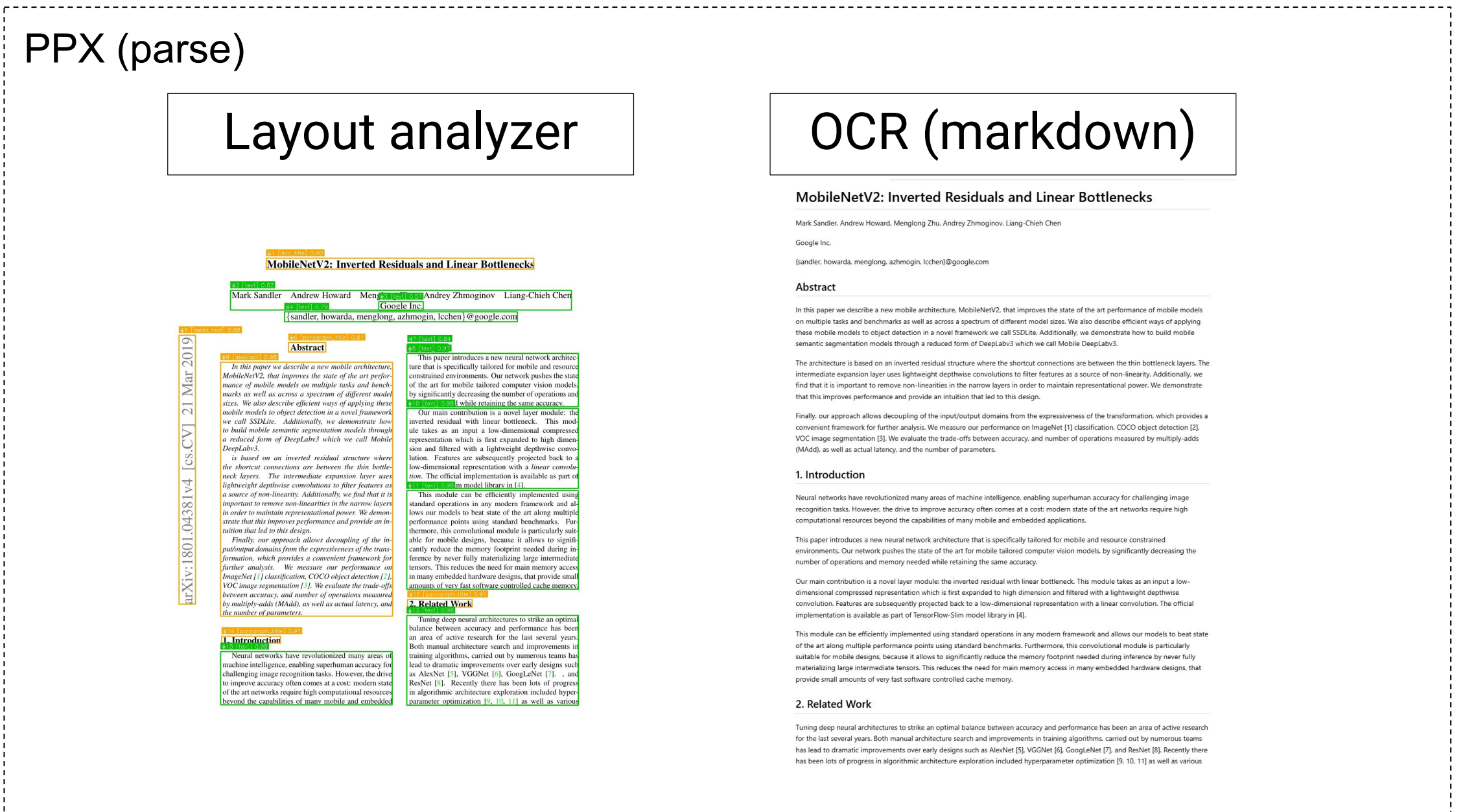
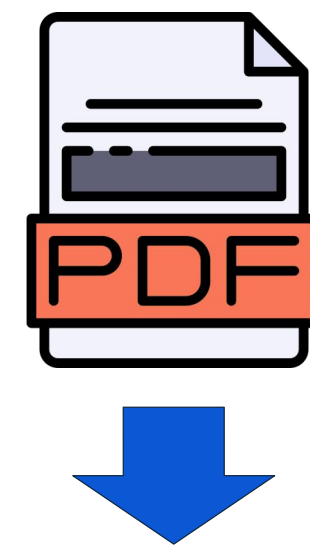
2. Methodology

- The system will consist of two (2) workflows: **document parsing** and **document search**.
- **Document Parsing:** Analyze layout, generate markdown using PaddleOCR and LLM tools (MistralOCR, Gemini, OpenAI). [1, 3]
- **Document Search:** Align markdown with document layout to establish *semantic provenance*. Generate embeddings using sentence transformers. [2]
- **Validation:** Measure alignment accuracy and verify returned source locations are correct.
- **Evaluation:** Measure Q&A accuracy, localization precision, and alignment accuracy on ground truth dataset.

4. References

- [1] C. Cui et al., "PaddleOCR 3.0 technical report," arXiv:2507.05595, 2025.
- [2] N. Reimers and I. Gurevych, "Sentence-BERT: Sentence embeddings using Siamese BERT-networks," EMNLP-IJCNLP, 2019, pp. 3982–3992.
- [3] Mistral AI, "Mistral OCR," 2025.
mistral.ai/news/mistral-ocr
- [4] Icons by Good Ware from Flaticon (flaticon.com)

3. System



5. Current Progress and Future Work

